| Unit – 1 |
| --- |
| **Storage technology** |
| **Unit-01/Lecture -01** |

**Introduction to information storage and management (ISM)**

Information storage and management is the only subject of its kind to fill the knowledge gap in understanding varied components of modern information storage infrastructure, including virtual environments. It provides comprehensive learning of storage technology, which will enable you to make more informed decisions in an increasingly complex it environment. Ism builds a strong understanding of underlying storage technologies and prepares you to learn advanced concepts, technologies, and products. You will learn about the architectures, features, and benefits of intelligent storage systems; storage networking technologies such as FC-SAN, ip-SAN, NAS, object-based and unified storage; business continuity solutions such as backup, replication, and archive; the increasingly critical area of information security; and the emerging field of cloud computing.

**Introduction to storage technology**

Storage systems are inevitable for modern day computing. All known computing platforms ranging from handheld devices to large super computers use storage systems for storing data temporarily or permanently. Beginning from punch card which stores a few bytes of data, storage systems have reached to multi terabytes of capacities in comparatively less space and power consumption. This tutorial is intended to give an introduction to storage systems and components to the reader.

Storage definition ,

here are a few definitions of storage when refers to computers.

- A device capable of storing data. The term usually refers to mass storage devices, such as disk and tape drives.

- In a computer, storage is the place where data is held in an electromagnetic or optical form for access by a computer processor. (whatis.com)

- Computer data storage; often called storage or memory refer to computer components, devices and recording media that retain digital data used for computing for some interval of time of these, i like the definition coined out by wikipedia.com. Likes and dislikes apart, in basic terms, computer storage can be defined as " device or media stores data for later retrieval". From the definition, we can see that the storage device possess two features namely "storage" and "retrieval". A storage facility without retrieval options seems to be of no use . A storage device may store application programs, databases, media files etc....

as we see in modern day computers, storage devices can be found in many forms. Storage devices can be classified based on many criterions. Of them, the very basic is as we learned in schools ie; primary storage and secondary storage. Storage devices can be further classified based on the memory technology that they use, based on its data volatility etc...

**Storage technologies  - Storage caching [Rgpv/dec2011(10)]**

Storage caching is used to buffer blocks of data in order to minimize the utilization of disks or storage arrays and to minimize the read / write latency for storage access. Especially for write intensive scenarios such as virtual desktops, write caching is very beneficial as it can keep the storage latency even during peak times at a low level.
Storage cache can be implemented in four places:

- disk (embedded memory – typically non-expansible)

- storage array (vendor specific embedded memory + expansion cards)

- computer accessing the storage (ram)

- storage network (i.e. Provisioning server)
- The cache can be subdivided into two categories:
- volatile cache: contained data is lost upon power outages (good for reads or non-critical writes)

- non-volatile cache: data is kept safe in case of power outages (good for reads and writes). Often referred as battery backed write cache .

To further increase the speed of the disk or storage array advanced algorithms such as read-ahead / read-behind or command queuing are commonly used.

| S.no | Rgpv questions | Year | Marks |
|------|----------------|------|-------|
| Q.1 | Explain storage technologies in detail? | Dec 2011 | 10 |

## Unit-01/Lecture-02

**Data proliferation – [Rgpv/ Dec 2015(2), Rgpv/dec2013(10), Rgpv/dec2013(5), Rgpv/dec2011(5)]**

Data proliferation refers to the prodigious amount of data, structured and unstructured, that businesses and governments continue to generate at an unprecedented rate and the usability problems that result from attempting to store and manage that data. While originally pertaining to problems associated with paper documentation, data proliferation has become a major problem in primary and secondary data storage on computers.

While digital storage has become cheaper, the associated costs, from raw power to maintenance and from metadata to search engines, have not kept up with the proliferation of data. Although the power required to maintain a unit of data has fallen, the cost of facilities which house the digital storage has tended to rise.

**Problems caused [Rgpv Dec2015(3), Rgpv Dec2014(2)]**

The problem of data proliferation is affecting all areas of commerce as the result of the availability of relatively inexpensive data storage devices. This has made it very easy to dump data into secondary storage immediately after its window of usability has passed. This masks problems that could gravely affect the profitability of businesses and the efficient functioning of health services, police and security forces, local and national governments, and many other types of organizations. Data proliferation is problematic for several reasons:

- Difficulty when trying to find and retrieve information. Increased manpower requirements to manage increasingly chaotic data storage resources.
- Slower networks and application performance due to excess traffic as users search and search again for the material they need.
- High cost in terms of the energy resources required to operate storage hardware.

**Proposed solutions**

- Applications that better utilize modern technology

- Reductions in duplicate data (especially as caused by data movement)

- Improvement of metadata structures

- Improvement of file and storage transfer structures

- User education and discipline.

- The implementation of information lifecycle management solutions to eliminate low-value information as early as possible before putting the rest into actively managed long-term storage in which it can be quickly and cheaply accessed.

| S.no | Rgpv questions | Year | Marks |
|------|----------------|------|-------|
| Q.1 | What do you mean by data proliferation ?Explain data proliferation process and major problem associated with data proliferation? | Dec 2015<br>Dec 2014<br>Dec 2013<br>Dec 2013 | 3<br>3<br>7<br>10 |
| Q.2 | Explain in brief data proliferation? | Dec 2015/<br>Dec 2011 | 2<br>5 |

# Unit-01/Lecture -03

**Overview of storage infrastructure components – [Rgpv/ dec 2015(7), Rgpv/dec2013 (7)]**

The choice of hard discs can have a profound impact on the capacity, performance and long-term reliability of any storage infrastructure. But it's unwise to trust valuable data to any single point of failure, so hard discs are combined into groups that can boost performance and offer redundancy in the event of disc faults. At an even higher level, those arrays must be integrated into the storage infrastructure -combining storage with network technologies to make data available to users over a lan or wan. If you're new to storage, or just looking to refresh some basic concepts, this chapter on data storage components can help to bring things into focus.

**The lowest level: hard discs**

Hard discs are random-access storage mechanisms that relegate data to spinning platters coated with extremely sensitive magnetic media. Magnetic read/write heads step across the radius of each platter in set increments, forming concentric circles of data dubbed "tracks." hard disc capacity is loosely defined by the quality of the magnetic media (bits per inch) and the number of tracks. Thus, a late-model drive with superior media and finer head control can achieve far more storage capacity than models just six-12 months old. Some of today's hard drives can deliver up to 750 gbytes of capacity. Capacity is also influenced by specific drive technologies including perpendicular recording, which fits more magnetic points into the same physical disc area.

**Grouping the discs: raid**

Hard discs are electromechanical devices and their working life is finite. Media faults, mechanical wear and electronic failures can all cause problems that render drive contents inaccessible. This is unacceptable for any organization, so tactics are often implemented to protect against failure. One of the most common data protection tactics is arranging groups of discs into arrays. This is known

as a raid.

Raid implementations typically offer two benefits; data redundancy and enhanced performance. Redundancy is achieved by copying data to two or more discs -when a fault occurs on one hard disc, duplicate data on another can be used instead. In many cases, file contents are also spanned (or striped) across multiple hard discs. This improves performance because the various parts of a file can be accessed on multiple discs simultaneously -rather than waiting for a complete file to be accessed from a single disc. Raid can be implemented in a variety of schemes, each with its own designation:

- Raid-0 -- disc striping is used to improve storage performance, but there is no redundancy.

- Raid-1 -- disc mirroring offers disc-to-disc redundancy, but capacity is reduced and performance is only marginally enhanced.

- Raid-5 -- parity information is spread throughout the disc group, improving read performance and allowing data for a failed drive to be reconstructed once the failed drive is replaced.

- Raid-6 -- multiple parity schemes are spread throughout the disc group, allowing data for up to two simultaneously failed drives to be reconstructed once the failed drive(s) are replaced.

There are additional levels, but these four are the most common and widely used. It is also possible to mix raid levels in order to obtain greater benefits. Combinations are typically denoted with two digits. For example, raid-50 is a combination of raid-5 and raid-0, sometimes noted as raid-5+0. As another example, raid-10 is actually raid-1 and raid-0 implemented together, raid-1+0. For more information on raid controllers, see the searchstorage.com article the new breed of raid controllers.

**Getting storage on the network**

Storage is useless unless network users can access it. There are two principle means of attaching

storage systems: NAS and SAN. NAS boxes are storage devices behind an ethernet interface, effectively connecting discs to the network through a single ip address. NAS deployments are typically straightforward and management is light, so new NAS devices can easily be added as more storage is needed. The downside to NAS is performance - storage traffic must compete for NAS access across the ethernet cable. But NAS access is often superior to disc access at a local server.

The SAN overcomes common server and NAS performance limitations by creating a sub network of storage devices interconnected through a switched fabric like FC or iscsi (called internet scsi or scsi-over-ip. Both FC and iscsi approaches make any storage device visible from any host, and offer much more availability for corporate data. FC is costlier, but offers optimum performance, while iscsi is cheaper, but somewhat slower. Consequently, FC is found in the enterprise and iscsi commonly appears in small and mid-sized businesses. However, SAN deployments are more costly to implement (in terms of switches, cabling and host bus adapters) and demand far more management effort.

| S.no | Rgpv questions | Year | Marks |
|------|----------------|------|-------|
| Q.1 | Explain briefly the evolution of the storage management? | Dec 2013 | 7 |
| Q.2 | Discuss Storage infrastructure components. | Dec 2015 | 7 |

## Unit-01/Lecture -04

**Information lifecycle management - [Rgpv/dec 2015(7),Rgpv/dec2014(7), Rgpv/dec 2013 (7), Rgpv/dec 2012 (5)]**
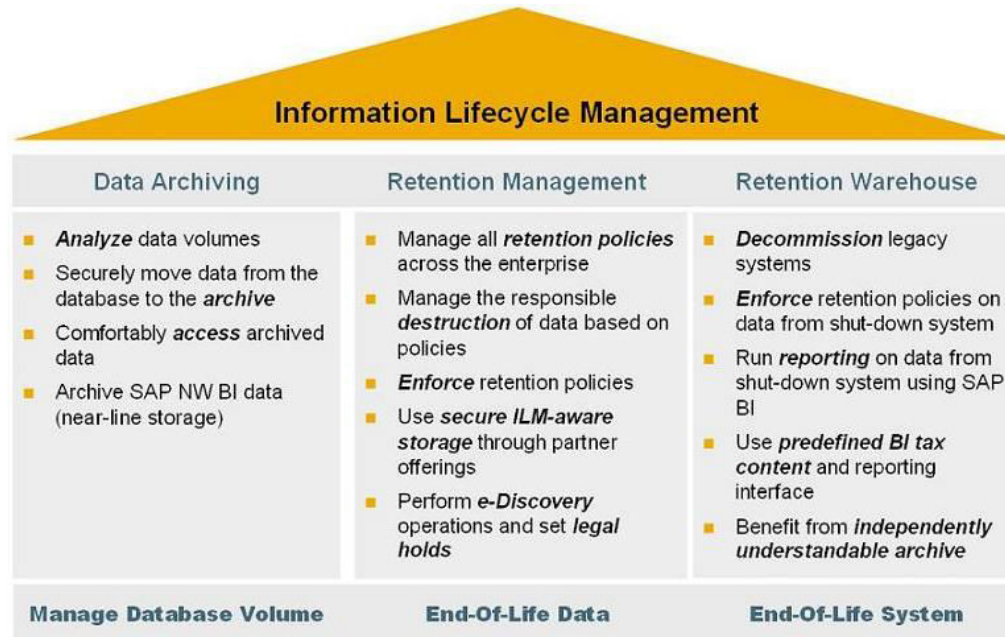
Information life cycle management (ILM) is a comprehensive approach to managing the flow of an information system's data and associated metadata from creation and initial storage to the time when it becomes obsolete and is deleted. Unlike earlier approaches to data storage management, ILM involves all aspects of dealing with data, starting with user practices, rather than just automating storage procedures, as for example, hierarchical storage management (HSM) does. Also in contrast to older systems, ilm enables more complex criteria for storage management than data age and frequency of access.

ILM products automate the processes involved, typically organizing data into separate tiers according to specified policies, and automating data migration from one tier to another based on those criteria. As a rule, newer data, and data that must be accessed more frequently, is stored on faster, but more expensive storage media, while less critical data is stored on cheaper, but slower media. However, the ILM approach recognizes that the importance of any data does not rely solely on its age or how often it's accessed. Users can specify different policies for data that declines in value at different rates or that retains its value throughout its life span. A path management application, either as a component of ILM software or working in conjunction with it, makes it possible to retrieve any data stored by keeping track of where everything is in the storage cycle.

ILM is often considered a more complex subset of data life cycle management (DLM).DLM products deal with general attributes of files, such as their type, size, and age; ILM products have more complex capabilities. For example, a DLM product would allow you to search stored data for a certain file type of a certain age. while an ilm product would let you search various types of stored files for instances of a specific piece of data, such as a customer number.

Data management has become increasingly important as businesses face compliance issues in the wake of legislation, that regulates how organizations must deal with particular types of data. Data management experts stress that information life cycle management should be an organization-wide enterprise, involving procedures and practices as well as applications.

"information lifecycle management comprises the policies, processes, practices, and tools used to align the business value of information with the most appropriate and cost effective it infrastructure from the time information is conceived through its final disposition. Information is aligned with business processes through management policies and service levels associated with applications, metadata, information, and data.

### Information Lifecycle Management

| Data Archiving | Retention Management | Retention Warehouse |
| --- | --- | --- |
| ■ *Analyze* data volumes<br>■ Securely move data from the database to the *archive*<br>■ Comfortably *access* archived data<br>■ Archive SAP NW BI data (near-line storage) | ■ Manage all *retention policies* across the enterprise<br>■ Manage the responsible *destruction* of data based on policies<br>■ *Enforce* retention policies<br>■ Use *secure ILM-aware storage* through partner offerings<br>■ Perform *e-Discovery* operations and set *legal holds* | ■ *Decommission* legacy systems<br>■ *Enforce* retention policies on data from shut-down system<br>■ Run *reporting* on data from shut-down system using SAP BI<br>■ Use *predefined BI tax content* and reporting interface<br>■ Benefit from *independently understandable archive* |
| **Manage Database Volume** | **End-Of-Life Data** | **End-Of-Life System** |

**Operations**

Operational aspects of ILM include backup and data protection; disaster recovery, restore, and restart; archiving and long-term retention; data replication; and day-to-day processes and procedures necessary to manage a storage architecture.

**Functionality**

For the purposes of business records, there are five phases identified as being part of the lifecycle continuum along with one exception. These are:

- Creation and receipt
- Distribution
- Use
- Maintenance
- Disposition

Creation and receipt deals with records from their point of origination. This could include their creation by a member of an organization at varying levels or receipt of information from an external source. It includes correspondence, forms, reports, drawings, computer input/output, or other sources.

Distribution is the process of managing the information once it has been created or received. This includes both internal and external distribution, as information that leaves an organization becomes a record of a transaction with others.

Use takes place after information is distributed internally, and can generate business decisions, document further actions, or serve other purposes.

Maintenance is the management of information. This can include processes such as filing, retrieval and transfers. While the connotation of 'filing' presumes the placing of information in a prescribed container and leaving it there, there is much more involved. Filing is actually the process of arranging information in a predetermined sequence and creating a system to manage it for its useful existence within an organization. Failure to establish a sound method for filing information makes its retrieval and use nearly impossible. Transferring information refers to the process of responding to requests, retrieval from files and providing access to users authorized by the organization to have access to the information. While removed from the files, the information is

tracked by the use of various processes to ensure it is returned and/or available to others who may need access to it.

Disposition is the practice of handling information that is less frequently accessed or has met its assigned retention periods. Less frequently accessed records may be considered for relocation to an 'inactive records facility' until they have met their assigned retention period. "although a small percentage of organizational information never loses its value, the value of most information tends to decline over time until it has no further value to anyone for any purpose. The value of nearly all business information is greatest soon after it is created and generally remains active for only a short time --one to three years or so-- after which its importance and usage declines. The record then makes its life cycle transition to a semi-active and finally to an inactive state." [1] retention periods are based on the creation of an organization-specific retention schedule, based on research of the regulatory, statutory and legal requirements for management of information for the industry in which the organization operates. Additional items to consider when establishing a retention period are any business needs that may exceed those requirements and consideration of the potential historic, intrinsic or enduring value of the information. If the information has met all of these needs and is no longer considered to be valuable, it should be disposed of by means appropriate for the content. This may include ensuring that others cannot obtain access to outdated or obsolete information as well as measures for protection privacy and confidentiality.'

Long-term records are those that are identified to have a continuing value to an organization. Based on the period assigned in the retention schedule, these may be held for periods of 25 years or longer, or may even be assigned a retention period of "indefinite" or "permanent". The term "permanent" is used much less frequently outside of the federal government, as it is not feasible to establish a requirement for such a retention period. There is a need to ensure records of a continuing value are managed using methods that ensure they remain persistently accessible for length of the time they are retained. While this is relatively easy to accomplish with paper or microfilm based records by providing appropriate environmental conditions and adequate protection from potential hazards, it is less simple for electronic format records. There are unique

concerns related to ensuring the format they are generated/captured in remains viable and the media they are stored on remains accessible. Media is subject to both degradation and obsolescence over its lifespan, and therefore, policies and procedures must be established for the periodic conversion and migration of information stored electronically to ensure it remains accessible for its required retention periods.

| S.no | Rgpv questions | Year | Marks |
|------|----------------|------|-------|
| Q.1 | What are the different phases of information life cycle model? | Dec 2015 | 7 |
| | | Dec 2014 | 7 |
| | | Dec 2013 | 7 |
| | | Dec 2011 | 10 |
| Q.2 | Explain briefly information life cycle implementation? | Dec 2012 | 5 |

## Unit-01/Lecture -05

**Data Categorization – [Rgpv/ dec 2015(2),Rgpv/dec2013(7) Rgpv/dec2013(10), Rgpv/dec2011(5)]]**

Data classification is the categorization of data for its most effective and efficient use. In a basic approach to storing computer data, data can be classified according to its critical value or how often it needs to be accessed, with the most critical or often-used data stored on the fastest media while other data can be stored on slower (and less expensive) media. This kind of classification tends to optimize the use of data storage for multiple purposes - technical, administrative, legal, and economic.

Data can be classified according to any criteria, not only relative importance or frequency of use. For example, data can be broken down according to its topical content, file type, operating platform, average file size in megabytes or gigabytes, when it was created, when it was last accessed or modified, which person or department last accessed or modified it, and which personnel or departments use it the most. A well-planned data classification system makes essential data easy to find. This can be of particular importance in risk management, legal discovery, and compliance with government regulations.

Computer programs exist that can help with data classification, but in the end it is a subjective business and is often best done as a collaborative task that considers business, technical, and other points-of-view.

**Data collections**

Data stewards may wish to assign a single classification to a collection of data that is common in purpose or function. When classifying a collection of data, the most restrictive classification of any of the individual data elements should be used. For example, if a data collection consists of a student's name, address and social security number, the data collection should be classified

as restricted even though the student's name and address may be considered public information.

**Why is it important?**

Data classification provides several benefits. It allows an organization to inventory its information assets. In many cases, information asset owners aren't aware of all of the different types of data they hold. It also allows central it to work with departments to develop specific security requirements that can be readily utilized.

N the field of data management, data classification as a part of information lifecycle management (ILM) process can be defined as a tool for categorization of data to enable/help organization to effectively answer following questions:

- What data types are available?
- Where are certain data located?
- What access levels are implemented?
- What protection level is implemented and does it adhere to compliance regulations?

When implemented it provides a bridge between it professionals and process or application owners. It staff is informed about the data value and on the other hand management (usually application owners) understands better to what segment of data centre has to be invested to keep operations running effectively. This can be of particular importance in risk management, legal discovery, and compliance with government regulations. Data classification is typically a manual process; however, there are many tools from different vendors that can help gather information about the data.

**How to start process of data classification?**

Note that this classification structure is written from a data management perspective and therefore has a focus for text and text convertible binary data sources. Images, video, and audio files are highly structured formats built for industry standard api's and do not readily fit within

the classification scheme outlined below.

First step is to evaluate and divide the various applications and data as follows:

- Relational or tabular data (around 15% of non audio/video data)

  - Generally describes proprietary data which can be accessible only through application or application programming interfaces (api)
  - Applications that produce structured data are usually database applications.
  - This type of data usually brings complex procedures of data evaluation and migration between the storage tiers.
  - To ensure adequate quality standards, the classification process has to be monitored by subject matter experts.

- Semi-structured or poly-structured data (all other non audio/video data that does not conform to a system or platform defined relational or tabular form).

  - Generally describes data files that have a dynamic or non-relational semantic structure (e.g. Documents,xml,json,device or system log output,sensor output).
  - Relatively simple process of data classification is criteria assignment.
  - Simple process of data migration between assigned segments of predefined storage tiers.

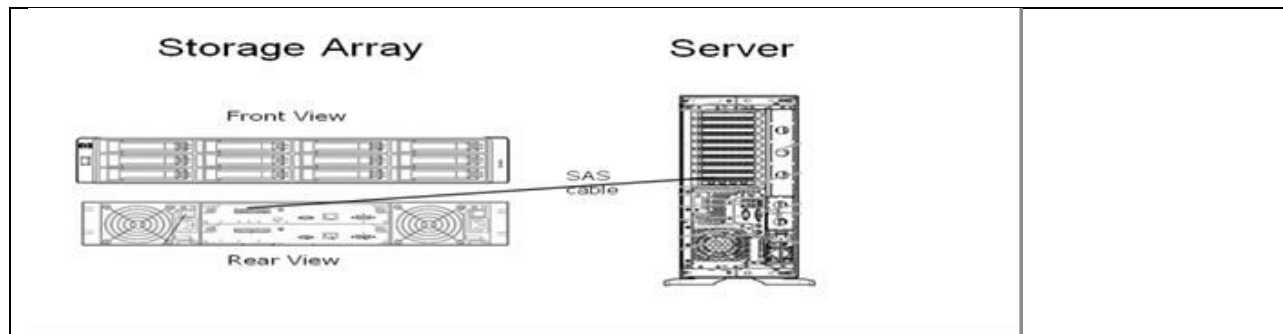| S.no | Rgpv question | Year | Marks |
|------|---------------|------|-------|
| Q.1 | What is data categorization ?why it is required. | Dec 2013 | 7 |
| Q.2 | What is data categorization? Explain challenges for data categorization | Dec 2013 | 10 |
| Q.3 | .<br>Explain briefly data categorization? | Dec 2015<br>Dec 2011 | 2<br>5 |

# Unit 01/Lecture - 06

**Evolution of various storage technologies[Rgpv/dec2012(10)],Rgpv/dec2011(5)]**

**Das (Direct Attached Storage)**

When windows servers leave the factory, they can be configured with several storage options. Most servers will contain 1 or more local disk drives which are installed internal to the server's cabinet. These drives are typically used to install the operating system and user applications. If additional storage is needed for user files or databases, it may be necessary to configure direct attached storage (das).

Das is well suited for a small-to-medium sized business where sufficient amounts of storage can be configured at a low startup cost. The das enclosure will be a separate adjacent cabinet that contains the additional disk drives. An internal pci-based raid controller is typically configured in the server to connect to the storage. The sas (serial attached scsi) technology is used to connect the disk arrays as illustrated in the following example.

As mentioned, one of the primary benefits of das storage is the lower startup cost to implement. Managing the storage array is done individually as the storage is dedicated to a particular server. On the downside, there is typically limited expansion capability with das, and limited cabling options (1 to 4 meter cables). Finally, because the raid controller is typically installed in the server, there is a potential single point of failure for the das solution.

**SAN (Storage Area Networks) [Rgpv Dec 2014(2)]**

With storage area networks (SAN), we typically see this solution used with medium-to-large size businesses, primarily due to the larger initial investment. Sans require an infrastructure consisting of SAN switches, disk controllers, hbas (host bus adapters) and fibre cables. Sans leverage external raid controllers and disk enclosures to provide high-speed storage for numerous potential servers.

The main benefit to a SAN-based storage solution is the ability to share the storage arrays to multiple servers. This allows you to configure the storage capacity as needed, usually by a dedicated SAN administrator. Higher levels of performance throughput are typical in a SAN environment, and data is highly available through redundant disk controllers and drives. The disadvantages include a much higher startup cost for sans, and they are inherently much more complex to manage. The following diagram illustrates a typical SAN environment.

**NAS (network attached storage)**

A third type of storage solution exists that is a hybrid option called network attached storage (NAS). This solution uses a dedicated server or "appliance" to serve the storage array. The storage can be commonly shared to multiple clients at the same time across the existing ethernet network. The main difference between NAS and das and SAN is that NAS servers utilize file level transfers, while das and SAN solutions use block level transfers which are more efficient.

NAS storage typically has a lower startup cost because the existing network can be used. This can be very attractive to small-to-medium size businesses. Most NAS models implement the storage arrays as iscsi targets that can be shared across the networks. Dedicated iscsi networks can also be configured to maximize the network throughput. The following diagram shows how a NAS configuration might look.

NAS Server                 Clients

Ethernet

| S.no | Rgpv question | Year | Marks |
|------|---------------|------|-------|
| Q.1 | Define SAN (Storage Area Network) | Dec 2014 | 2 |
| Q.2 | Explain briefly about the evolution of storage technologies and architecture? | Dec 2012 | 10 |
| Q.3 | Explain briefly Evolution of various storage technologies | Dec 2011 | 5 |

## Unit 01/Lecture - 07

**Data Centre - [Rgpv/dec 2014(2), Rgpv/dec2013(7), Rgpv/dec2013(10), Rgpv/dec2012(10), Rgpv/dec2011(10)]**

A data center (sometimes spelled datacenter) is a centralized repository, either physical or virtual, for the storage, management, and dissemination of data and information organized around a particular body of knowledge or pertaining to a particular business.

The national climatic data center (ncdc), for example, is a public data center that maintains the world's largest archive of weather information. A private data center may exist within an organization's facilities or may be maintained as a specialized facility. Every organization has a data center, although it might be referred to as a server room or even a computer closet.

In that sense, data center may be synonymous with network operations center (noc), a restricted access area containing automated systems that constantly monitor server activity, web traffic, and network performance.

Organizations maintain data centers to provide centralized data processing capabilities across the enterprise. Data centers store and manage large amounts of mission-critical data. The data center infrastructure includes computers, stor- age systems, network devices, dedicated power backups, and environmental controls (such as air conditioning and fire suppression).

Large organizations often maintain more than one data center to distribute data processing workloads and provide backups in the event of a disaster. The storage requirements of a data center are met by a combination of various stor- age architectures.

**Core elements**

Five core elements are essential for the basic functionality of a data center:

- application: an application is a computer program that provides the logic for computing operations. Applications, such as an order processing system, can be layered on a database, which in turn uses operating system services to perform read/write operations to storage devices.
- Database: more commonly, a database management system (DBMS) provides a structured way to store data in logically organized tables that are interrelated. A dbms optimizes the storage and retrieval of data.
- server and operating system: a computing platform that runs applications and databases.
- network: a data path that facilitates communication between clients and servers or between servers and storage.
- storage array: a device that stores data persistently for subsequent use.

These core elements are typically viewed and managed as separate entities, but all the elements must work together to address data processing requirements.

**Key requirements for data center elements**

Uninterrupted operation of data centers is critical to the survival and success of a business. It is necessary to have a reliable infrastructure that ensures data is accessible at all times. While the requirements, shown in figure 1-6, are appli- cable to all elements of the data center infrastructure, our focus here is on storage systems.

- Availability: all data center elements should be designed to ensure acces- sibility. The inability of users to access data can have a significant negative impact on a business.
- Security: polices, procedures, and proper integration of the data cen- ter core elements that will prevent unauthorized access to information must be established.

In addition to the security measures for client access, specific mechanisms must enable servers to access only their allocated resources on storage arrays.

- Scalability: data center operations should be able to allocate additional processing capabilities or storage on demand, without interrupting busi- ness operations. Business growth often requires deploying more servers, new applications, and additional databases. The storage solution should be able to grow with the business.

- Performance: all the core elements of the data center should be able to provide optimal performance and service all processing requests at high speed. The infrastructure should be able to support performance requirements.

- Data integrity: data integrity refers to mechanisms such as error correc- tion codes or parity bits which ensure that data is written to disk exactly as it was received. Any variation in data during its retrieval implies cor- ruption, which may affect the operations of the organization.

- Capacity: data center operations require adequate resources to store and process large amounts of data efficiently. When capacity requirements increase, the data center must be able to provide additional capacity with- out interrupting availability, or, at the very least, with minimal disruption. Capacity may be managed by reallocation of existing resources, rather than by adding new resources.

- Manageability: a data center should perform all operations and activi- ties in the most efficient manner. Manageability can be achieved through automation and the reduction of human (manual) intervention in com- mon tasks.

① A customer places an order through the AUI of the order processing application software located on the client computer.

② The client connects to the server over the LAN and accesses the DBMS located on the server to update the relevant information such as the customer name, address, payment method, products ordered, and quantity ordered.

③ The DBMS uses the server operating system to read and write this data to the database located on physical disks in the storage array.

④ The Storage Network provides the communication link between the server and the storage array and transports the read or write commands between them.

⑤ The storage array, after receiving the read or write commands from the server, performs the necessary operations to store the data on physical disks.

| S.no. | Year | Year | Mark |
|-------|------|------|------|
| Q.1 | What are the data centre? What are the requirement for the design of a secure data centre. | Dec2013 Dec2012 Dec2012 | 7 10 10 |
| Q.2 | What is the significance of Data Center in storage technology | Dec 2014 | 2 |

**Reference**

| Book | Author | Priority |
|---|---|---|
| Information storage management | G. Somasundaram Alok Shrivastava | 1 |
| Storage Network explained : Basic and application of fiber channels, SAN, NAS, iSESI | Ulf Troppens, Wolfgang Mueller-Friedt, Rainer Erkens, Rainer Wolafka, Nils Haustein | 2 |
| Nick Antonopoulos, Lee Gillam | Cloud Computing : Principles, System & Application | 3 |